

# Improving Model Performance, Portability and Productivity with Apache TVM and the OctoML platform

*ModSim'21*

*Luis Ceze*

Co-founder & CEO, OctoML.

Professor, University of Washington.

*Credit goes to many collaborators at OctoML, UW, the Apache TVM community!*

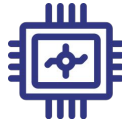
# The Problem: Machine Learning is hard and costly to deploy



Engineering



Cloud costs



Silicon resources



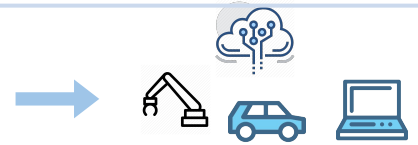
Data Collection and Processing



Model Development & Training



Model Optimization & Deployment



Cloud, Mobile, Edge Inference

Large gap not covered by typical "MLOps"



## Performance/efficiency gap

Model optimization and tuning is essential for viable deployment.



## Productivity gap

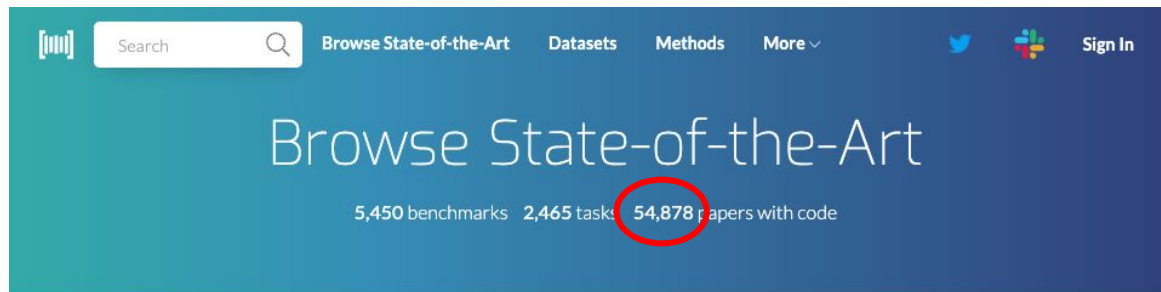
Weeks to months of effort to get a model ready for deployment.



## Portability gap

Changing deployment HW requires significant manual effort. Vendor lock-in.

# Trend: Machine learning workload diversity is exploding



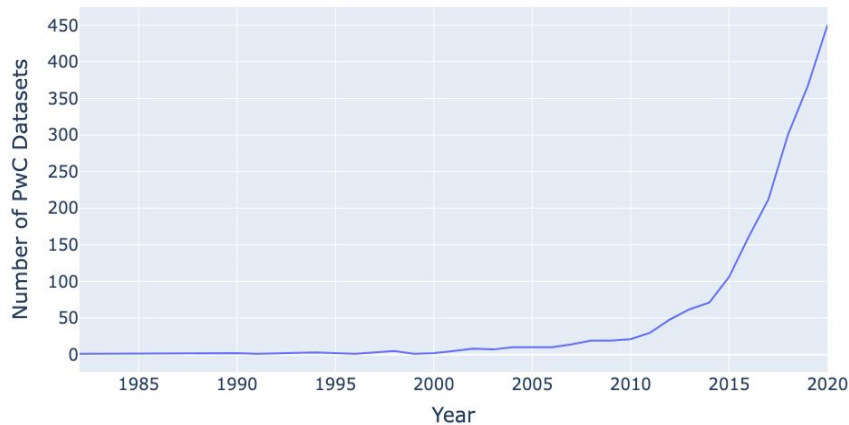
## Computer Vision

Task	Benchmarks	Papers with code
Semantic Segmentation	188	2099
Image Classification	260	1821
Object Detection	237	1576
Image Generation	160	698
Denosing	100	672

▶ See all 1129 tasks

# Trend: Machine learning workload diversity is exploding

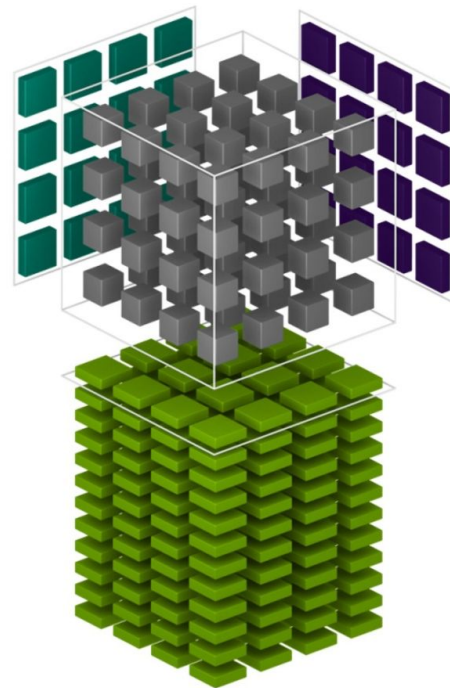
paperswithcode.com



# Trend: ML hardware capabilities exploding

E.g., TensorCores

- FP16
- Int8
- Int4, Int2, Int1
- Bfloat16, TF32
- FP64
- Tensor Core sparsity



$$\mathbf{D} = \begin{pmatrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{pmatrix} \begin{pmatrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{pmatrix} + \begin{pmatrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{pmatrix}$$

FP16 or FP32                      FP16                      FP16 or FP32

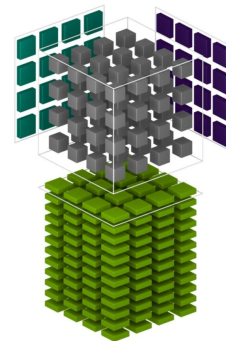
# Trend: ML hardware capabilities exploding

More programmability! E.g.,

- Asynchronous copy instruction loads data directly from global memory into shared memory, optionally bypassing L1 cache
- New instructions for L2 cache management and residency controls.
- New warp-level reduction instructions supported by CUDA Cooperative Groups.

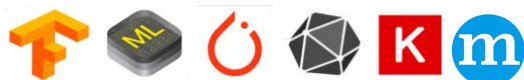
$$\mathbf{D} = \begin{pmatrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{pmatrix} \begin{pmatrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{pmatrix} + \begin{pmatrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{pmatrix}$$

FP16 or FP32                      FP16                      FP16 or FP32

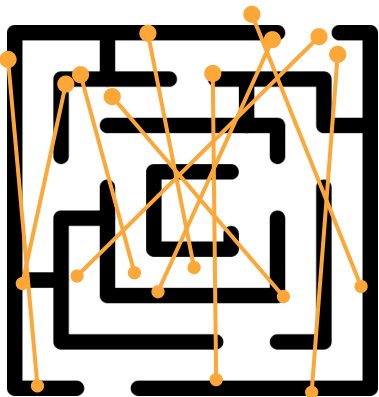


# An exploding ecosystem...

Rapidly evolving ML software ecosystem



- Use-case or HW-specific stack often hand-written
- Painful and unproductive for users
- Unsustainable for HW/platform vendors: need to keep up with model and framework evolution



Cambrian explosion of HW backends

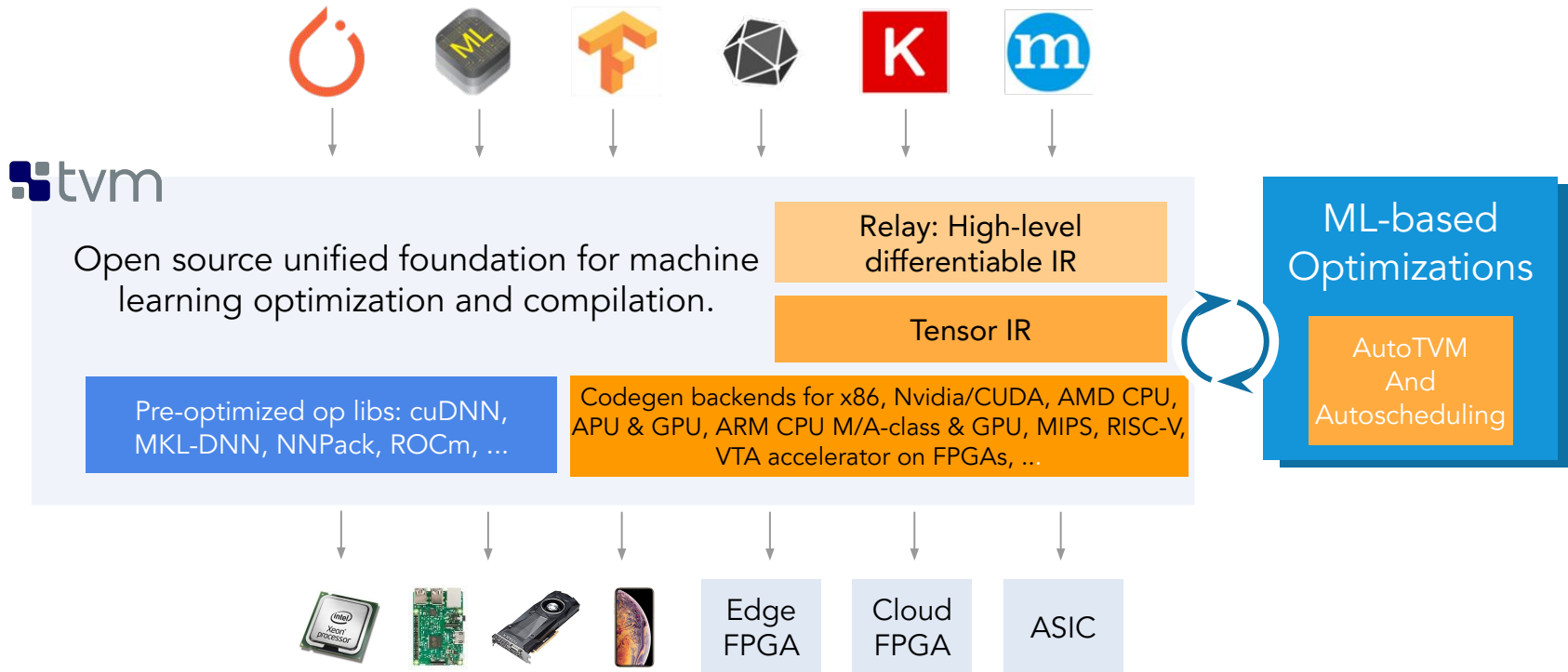


Automated, [open source](#), unified optimization and compilation framework for deep learning.

Model In, HW-specific, native code out.



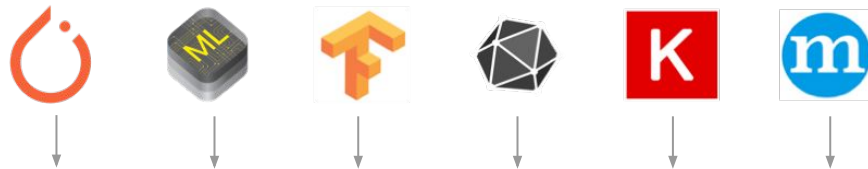
# TVM: Automatic ML optimizer, compiler and runtime





# TVM: Automatic ML optimizer, compiler and runtime

w/ state-of-the-art performance



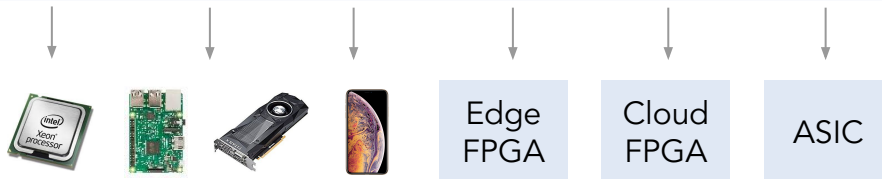
Open source unified foundation for machine learning optimization and compilation.

Relay: High-level differentiable IR

Tensor IR

Pre-optimized op libs: cuDNN, MKL-DNN, NNpack, ROCm, ...

Codegen backends for x86, Nvidia/CUDA, AMD CPU, APU & GPU, ARM CPU M/A-class & GPU, MIPS, RISC-V, VTA accelerator on FPGAs, ...



Reduce model time-to-market



Build your model once, run anywhere



Cut capital and operational ML costs

# TVM is an industry standard open source ML stack



Every "Alexa" wake-up today across all devices uses a model optimized with TVM



"[TVM enabled] real-time on mobile CPUs for free...We are excited about the performance TVM achieves." More than 85x speed-up for speech recognition model.



Bing query understanding: 112ms (Tensorflow) -> 34ms (TVM). QnA bot: 73ms->28ms (CPU), 10.1ms->5.5ms (GPU) - TVMconf 2019.



"TVM is key to ML Access on Hexagon" - Jeff Gehlhaar, VP Technology



Unified ML stack for CPU, GPU, NPU built on TVM. TVMconf 2020.



UNTETHER AI

Cross-product of {Models} x {Hardware} is large.  
Strong community support makes diversity manageable and ensures future-proofness.



Open source  
~570+ contributors from industry and academia.



13% HW vendors,  
35% research, 50%+  
ML end user.

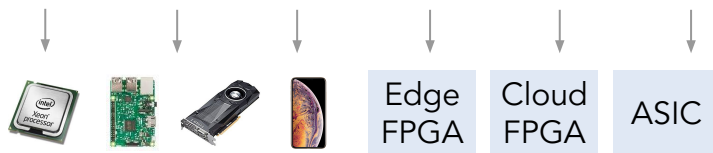
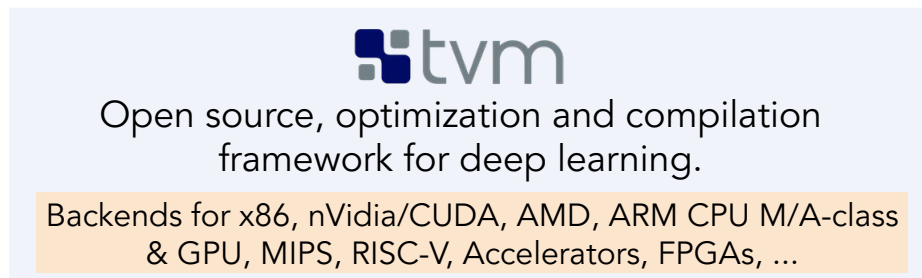


# Why use Apache TVM?



ML end user: *wants to quickly optimize their model on best/chosen HW target.*

- ML engineer for production apps
- Product R&D
- ML researchers

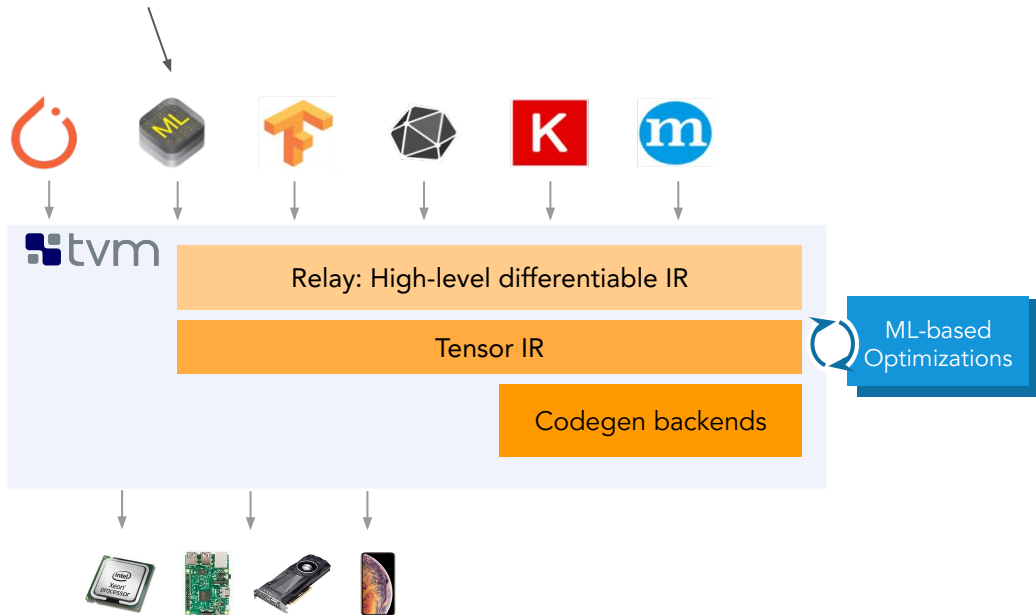


HW chip vendor & platform provider: *wants to offer their customers the best ML SW stack for all models*

- System SW engineers
- Technical sales demos

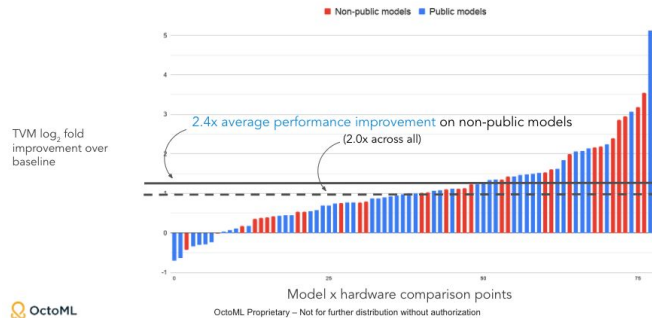
# TVM as a compiler and runtime framework

Need: End to end DL model optimization



- OctoML historical data
  - 2x avg improvement (up to 30x)

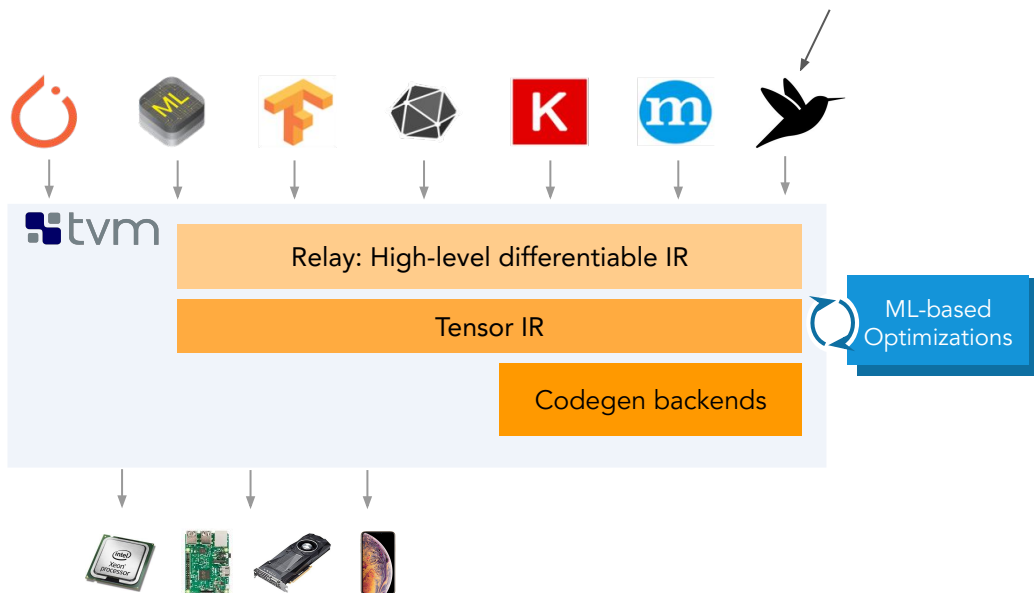
TVM Performance at OctoML



# TVM as a compiler and runtime framework

Need: End to end DL model optimization

Need: End to end classical ML model optimization



Case study:

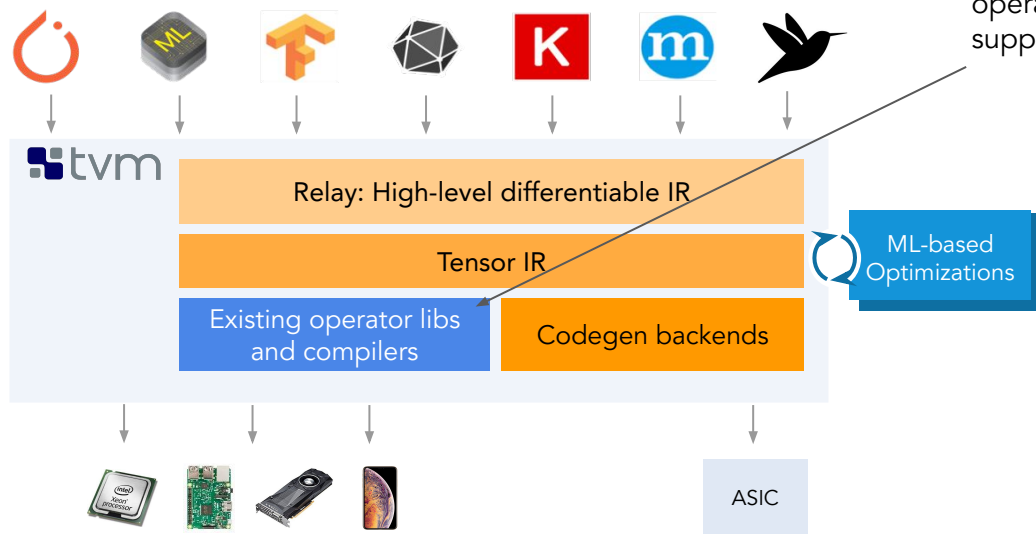
- Support decision tree optimization
  - Convert to sparse tensors, work w/ Microsoft
- Hummingbird
  - Microsoft Azure Data
  - 2-3x throughput increases
  - Broader algorithm support on GPU and future accelerators

# TVM as a compiler and runtime framework

Need: End to end DL model optimization

Need: End to end classical ML model optimization

Need: Enabling existing kernel operator libraries to support frameworks.



Case study:

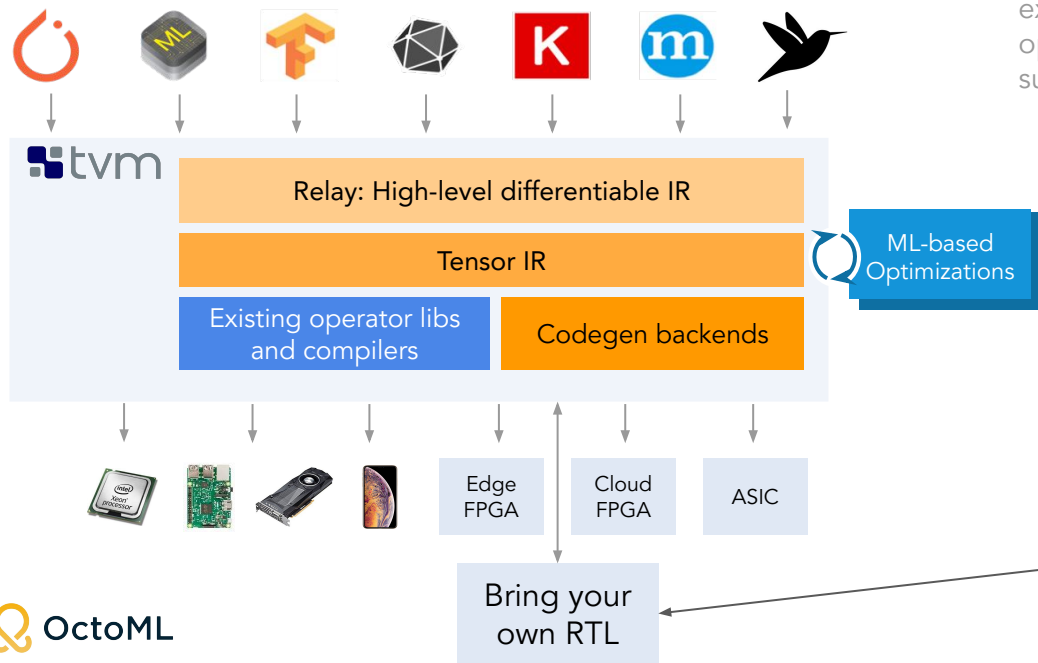
- ARM Ethos NPU
  - Both codegen and existing
  - Enables single compiler stack across CPU, GPU, and NPU
- Amazon Inferentia

# TVM as a compiler and runtime framework

Need: End to end DL model optimization

Need: End to end classical ML model optimization

Need: Enabling existing kernel operator libraries to support frameworks.



Case study:  
VTA Open Source DL accelerator and "Bring your own RTL"

- Enables rapid feedback loops for SW/HW codesign
- Verilator and other EDA simulators supported

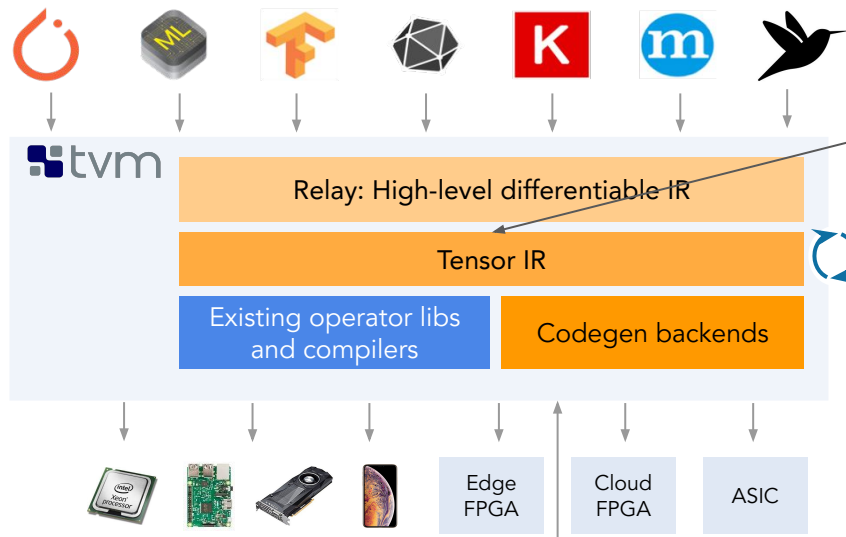
Need: Support rapid HW/SW codesign.

# TVM as a compiler and runtime framework

Need: End to end DL model optimization

Need: End to end classical ML model optimization

Need: Enabling existing kernel operator libraries to support frameworks.



Need: Enable users to write custom kernels with:

1. Full control (eg: a "nicer" CUDA)
2. Autotuning
3. Fully automatic autoscheduling

Case study:

Huggingface PruneBERT

- 2x acceleration on AMD CPUs
- 3x on NVIDIA and AMD GPUs

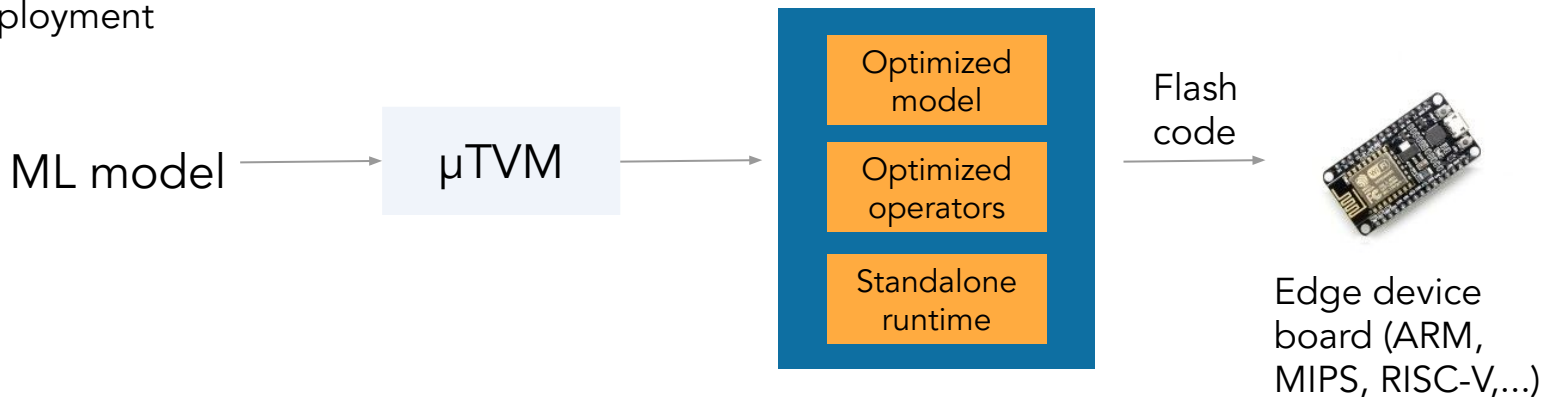
Kernel wise: 3-10x faster than cuBLAS and cuSPARSE

Single bit end to end acceleration (6-15x)



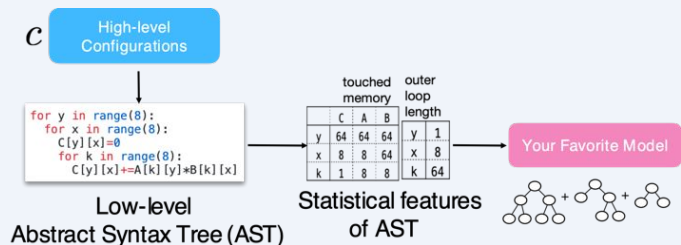
# $\mu$ TVM - Bare-metal model deployment for edge devices

Optimize, compile and package model for standalone bare metal deployment

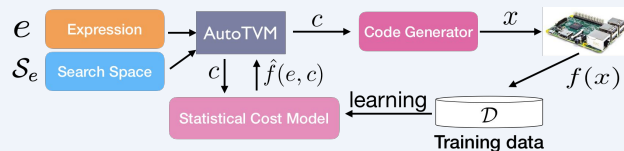


See recent demo on TVM for Azure Sphere deployment.

# ML-based optimizations



Benefit: low-level AST is a common representation (general, task invariant)

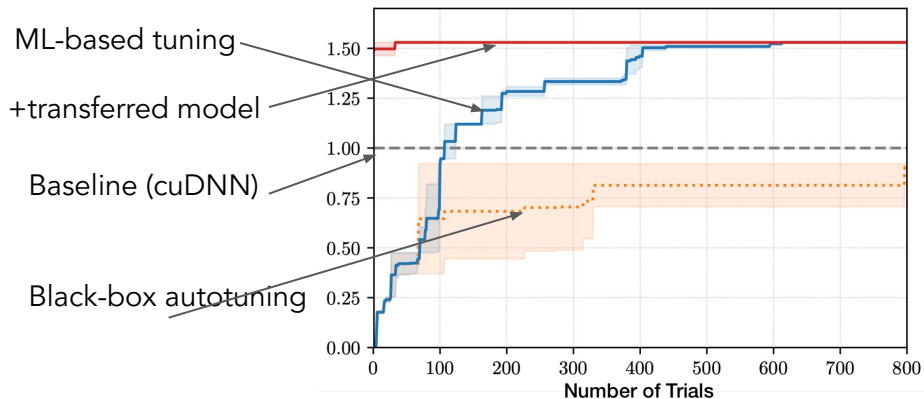


Automatically adapt to hardware type by learning, transferable to other tasks.

Extract hierarchical optimization search space from naive implementation (auto-scheduling).

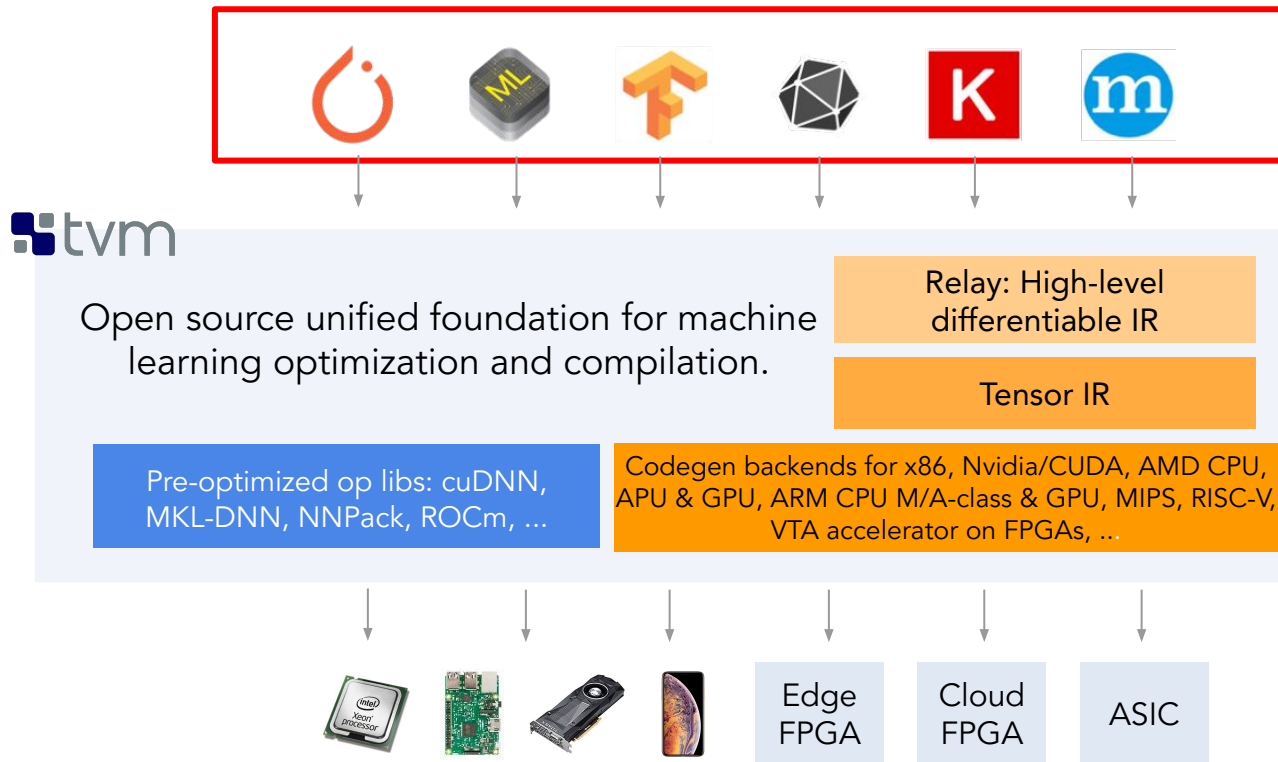
Start from repository of existing cost models, augment with more experiments (ML-based cost models).

The more it is used, the better it gets!





# TVM deep dive



Importers from a variety of formats into TVM.

- Caffe {2}
- CoreML
- Darknet
- Keras
- MXNet
- ONNX
- PyTorch
- TensorFlow
- TFLite

## Shape Information and Model Definition

```
→ shape_dict = {"input": x.shape}
   mod, params = relay.frontend.from_onnx(onnx_model, shape_dict)

   with tvn.transform.PassContext(opt_level=1):
       intrp = relay.create_executor("graph", mod, tvn.cpu(0), "llvm")

   tvn_output = intrp.evaluate()(tvn.nd.array(x), **params).asnumpy()
```

## Compile and Create Executor

```
shape_dict = {"input": x.shape}
mod, params = relay.frontend.from_onnx(onnx_model, shape_dict)

with tvm.transform.PassContext(opt_level=1):
    intrp = relay.create_executor("graph", mod, tvm.cpu(0), "llvm")

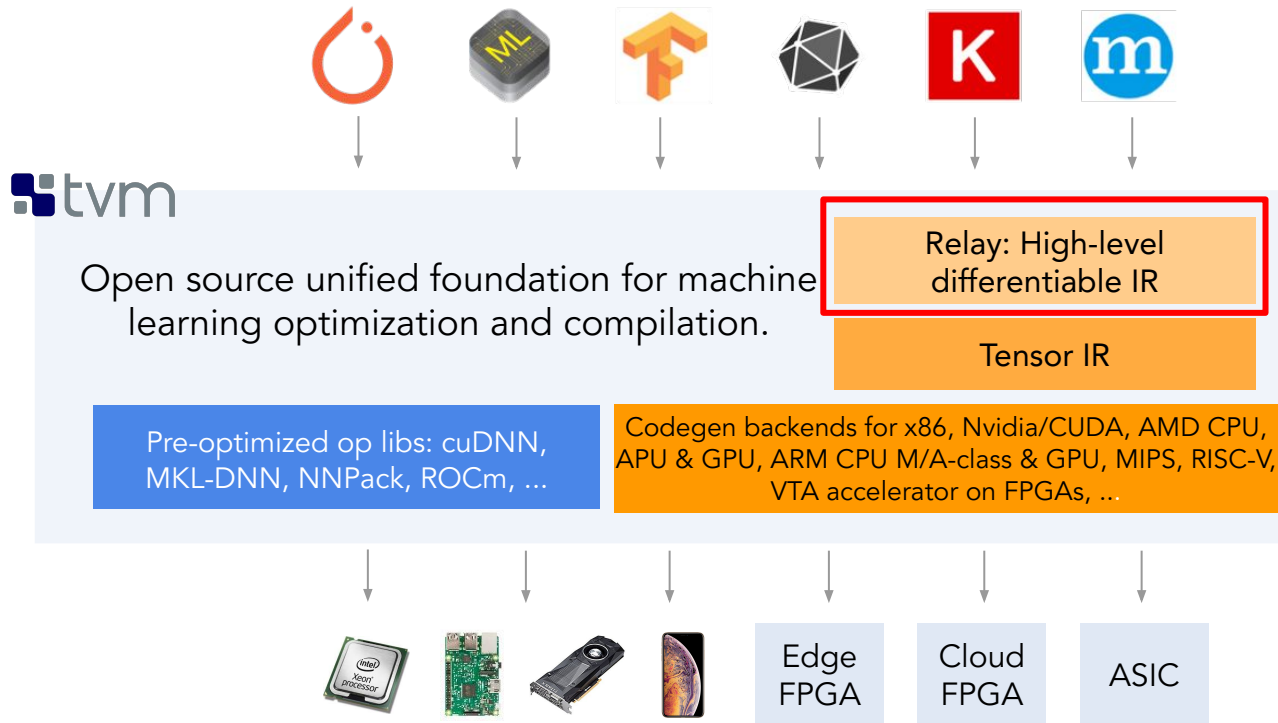
tvm_output = intrp.evaluate()(tvm.nd.array(x), **params).asnumpy()
```

```
shape_dict = {"input": x.shape}
mod, params = relay.frontend.from_onnx(onnx_model, shape_dict)

with tvm.transform.PassContext(opt_level=1):
    intrp = relay.create_executor("graph", mod, tvm.cpu(0), "llvm")

Run → tvm_output = intrp.evaluate()(tvm.nd.array(x), **params).asnumpy()
```

# TVM deep dive

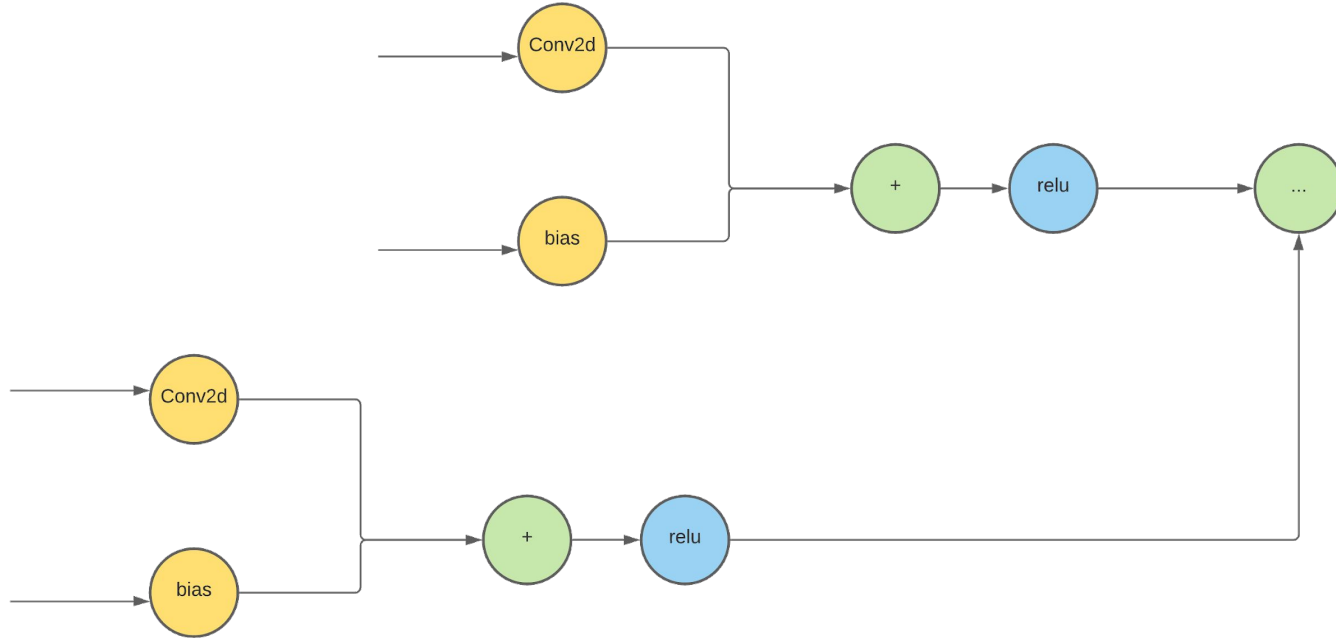


A functional graph-level IR, supports richer programming model.

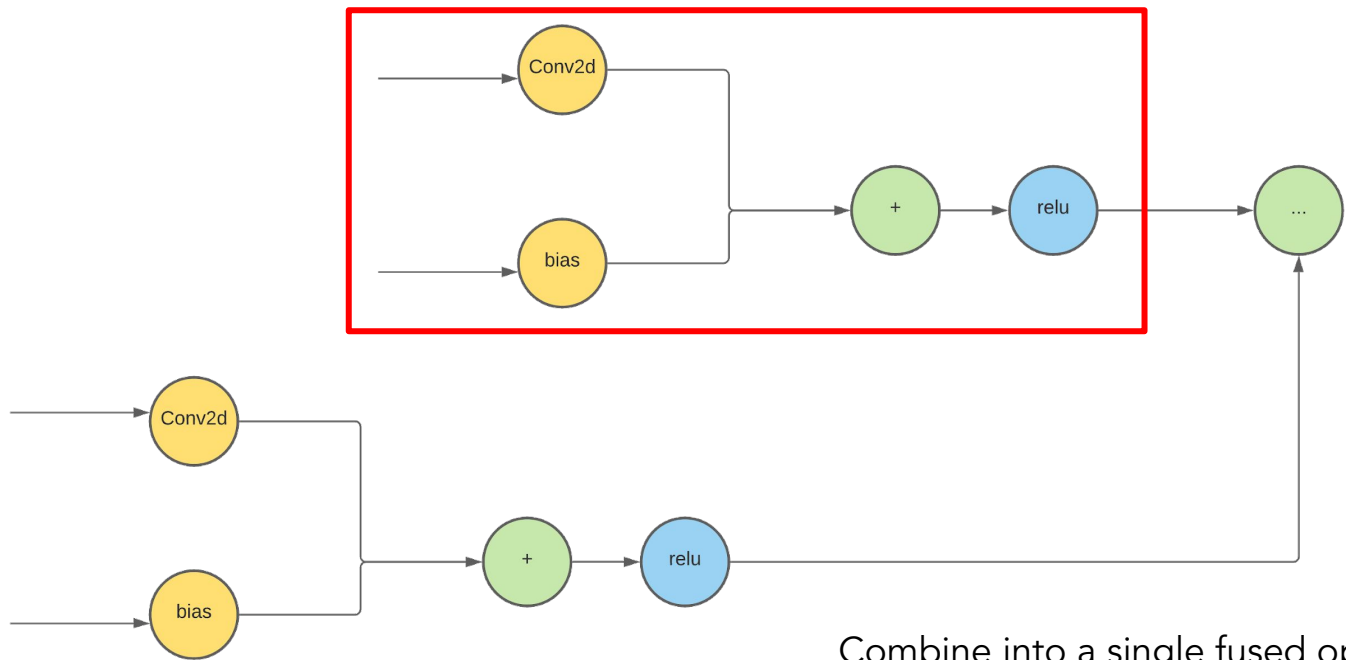
Functions, data types, primitive operations, tensors, and more.



# Relay

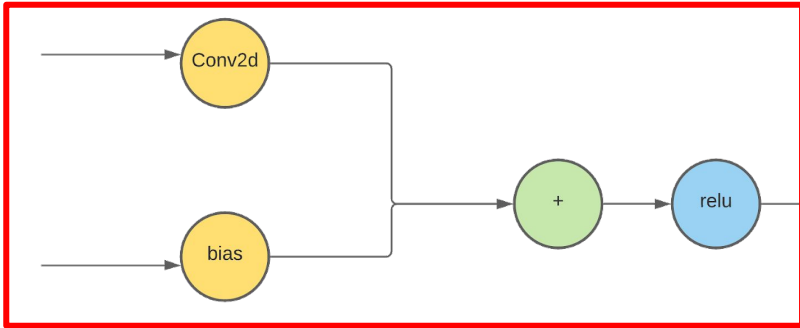
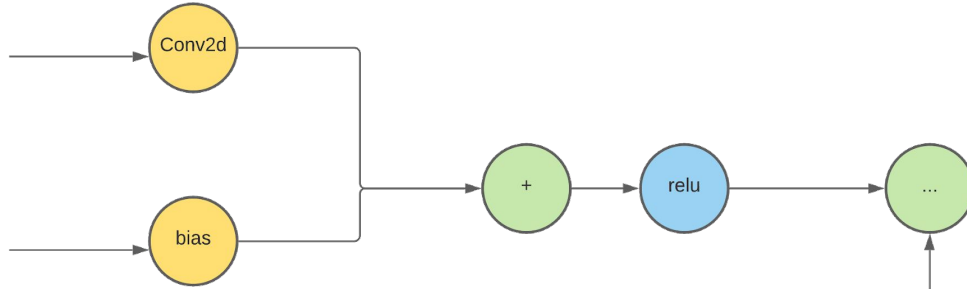


# Relay: Fusion



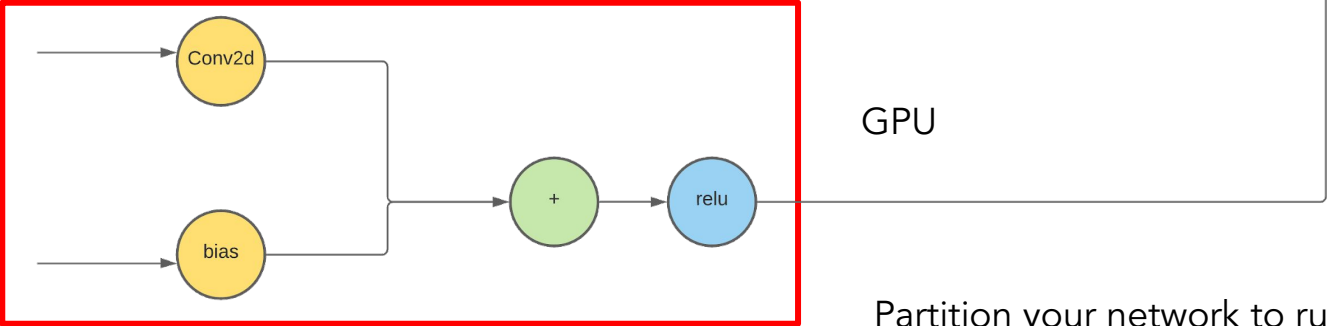
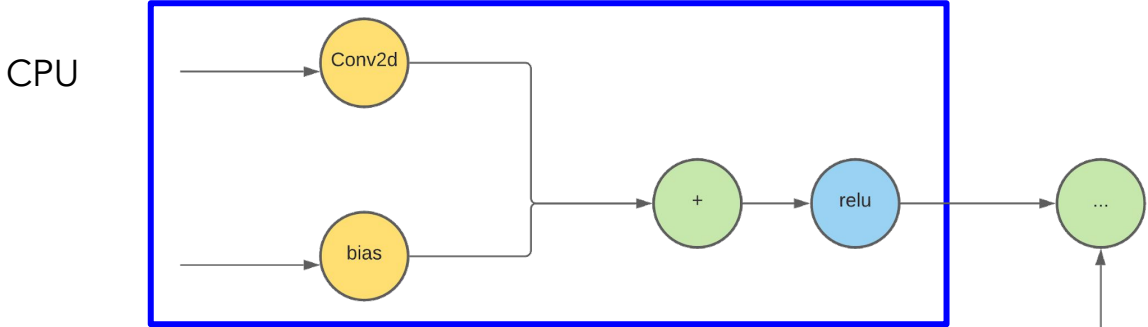
Combine into a single fused operation which can then be optimized specifically for your target.

# Relay: Fusion



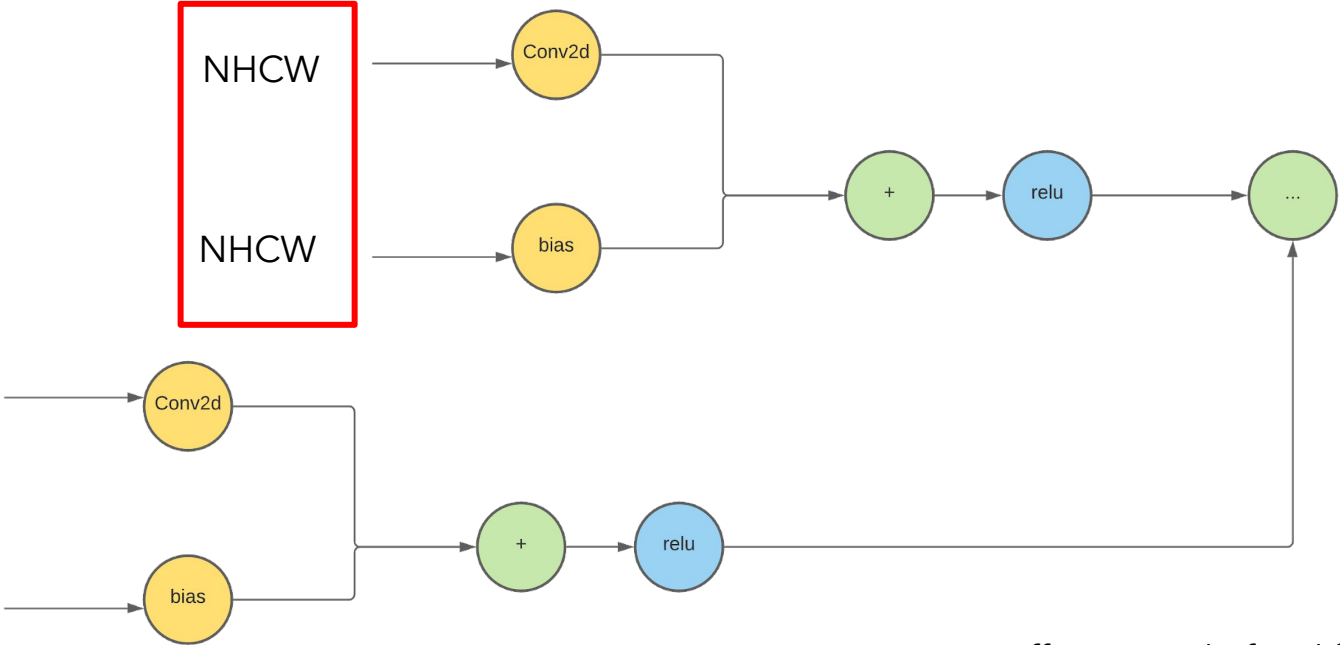
Combine into a single fused operation which can then be optimized specifically for your target.

# Relay: Device Placement



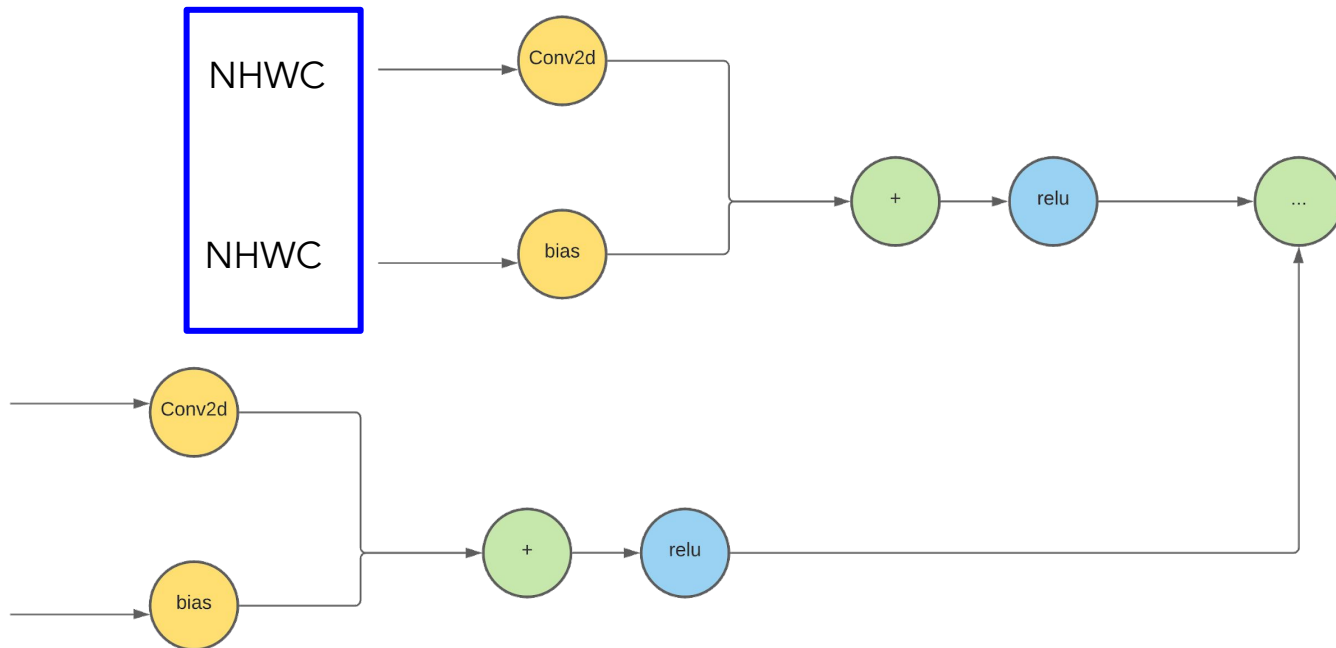
Partition your network to run on multiple devices.

# Relay: Layout Transformation



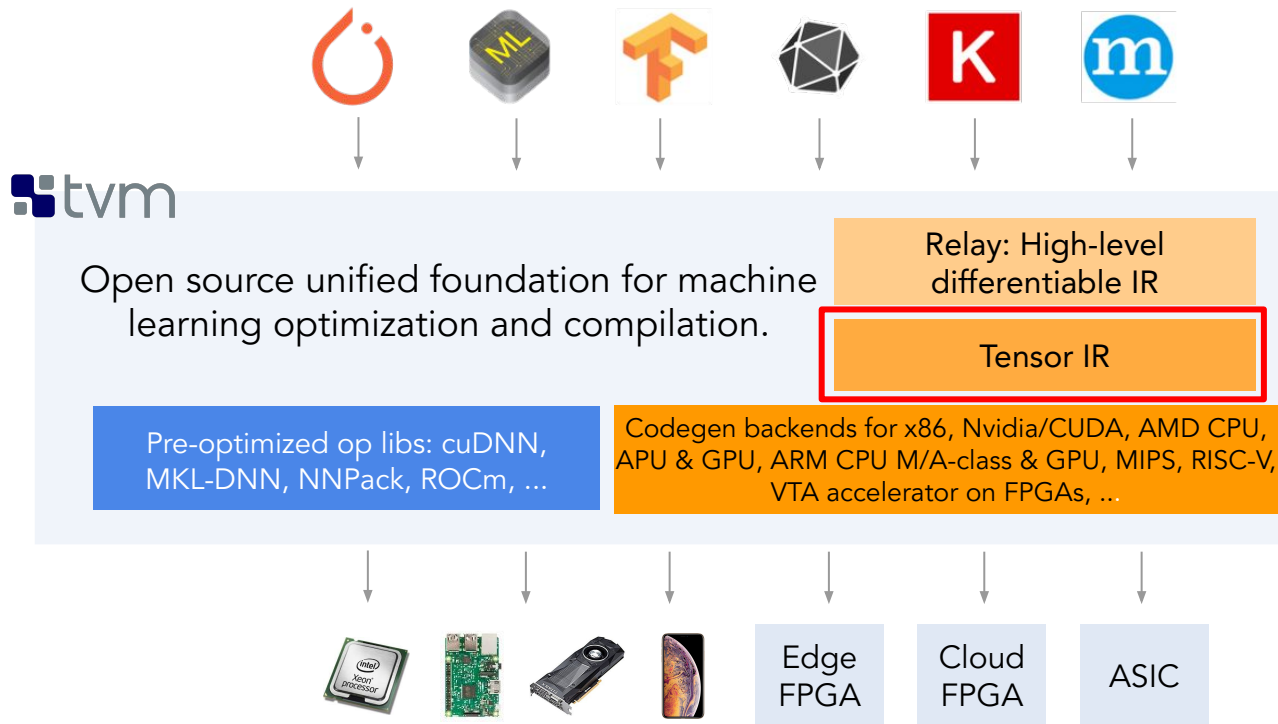
Generate efficient code for different data layouts.

# Relay: Layout Transformation



Generate efficient code for different data layouts.

# TVM deep dive



Low level IR represents kernels which can be optimized and compiled for all supported platforms.

# Tensor IR: The Kernel Enabler

```
@tvm.script.tir
class Module:
    def mmult(A: ty.handle, B: ty.handle, C: ty.handle) ->None:
        # function attr dict
        tir.func_attr({"global_symbol": "mmult", "tir.noalias": True})
        A_1 = tir.buffer_bind(A, [024, 024], ...)
        B_1 = tir.buffer_bind(B, [024, 024], ...)
        C_1 = tir.buffer_bind(C, [024, 024], ...)
        # body
        tir.attr(C_1, "realize_scope", "")
        tir.realize(C_1[0:1024, 0:1024])
        for x in tir.range(0, 1024):
            for y in tir.range(0, 1024):
                C_1[x, y] = tir.float32()
                for k in tir.range(0, 1024):
                    C_1[x, y] = (C_1[x, y] + (A_1[x, k]*B_1[k, y]))
```

- An imperative IR for describing loop nests and low level code.
- Can represent loop-y computations, perform allocation, bind tensors to backing buffers, and so on.
- Platform agnostic CUDA-like IR.
- Provides multiple target support via LLVM, and source code generation.
- Simple cross-platform runtime API which is implemented for CUDA, Vulkan, OpenCL, Metal and so on.



# Automating Kernel Code Generation is Key

## AutoTVM

Heavily template-based

Tensorization (limited)

Small design space

Customize design space

Slower compilation

## AutoScheduling (e.g., Anso)

Template-free

No tensorization

Large design space

Uncustomizable

Faster compilation



Automatic tensorization?

Perf on quantized model?

Training (AutoDiff)?

Dynamic shape?

Faster tuning?



# Auto Tensorization

- Tensor intrinsics are important
  - [NVIDIA Tensor Core](#)
  - Intel VNNI
  - ARM dot
  - ...
- Tensorization is hard in Anso
  - Handling structural matching / rewriting?
- Auto Tensorization in meta schedule
  - Just a search rule!

# Meta Schedule: Auto Tensorization in 4 Steps!

- Step 1: Describe the Tensor Intrinsic!

```
@tvm.script.tir
def tensorcore_desc(a: ty.handle, b: ty.handle, c: ty.handle):
    ...
    for i, j, k in tir.grid(16, 16, 16):
        ...
        C[vii, vjj] = C[vii, vjj] + A[vii, vkk] * B[vjj, vkk]

@tvm.script.tir
def tensorcore_impl(a: ty.handle, b: ty.handle, c: ty.handle):
    ...
    tir.evaluate(tir.tvm_mma_sync(
        C.data, C.elem_offset // 256,
        A.data, A.elem_offset // 256,
        B.data, B.elem_offset // 256,
        C.data, C.elem_offset // 256,
        dtype="handle",
    ))
```

# Meta Schedule: Auto Tensorization in 4 Steps!

- Step 2: Automatic structural fuzzy match

```
@tvm.script.tir
def batch_matmul(a: ty.handle, b: ty.handle, c: ty.handle):
    ...
    for n, i, j, k in tir.grid(16, 128, 128, 128):
        ...
        C[vn, vi, vj] += A[vn, vi, vk] * B[vn, vj, vk]
```

```
@tvm.script.tir
def tensorcore_desc(a: ty.handle, b: ty.handle, c: ty.handle):
    ...
    for i, j, k in tir.grid(16, 16, 16):
        ...
        C[vii, vjj] = C[vii, vjj] + A[vii, vkk] * B[vjj, vkk]
```



# Meta Schedule: Auto Tensorization in 4 Steps!

- Step 3: Automatic loop re-structuring & Mark tensorize region

```
@tvm.script.tir
def batch_matmul(a: ty.handle, b: ty.handle, c: ty.handle):
    ...
    for i0, i1_o, i2_o, i3_o in tir.grid(16, 8, 8, 8):
        ...
        for i1_i, i2_i, i3_i in tir.grid(16, 16, 16):
            ...
            C[vn_1, vi_1, vj_1] += A[vn_1, vi_1, vk_1] * B[vn_1, vj_1, vk_1]
```

```
@tvm.script.tir
def tensorcore_desc(a: ty.handle, b: ty.handle, c: ty.handle):
    ...
    for i, j, k in tir.grid(16, 16, 16):
        ...
        C[vii, vjj] = C[vii, vjj] + A[vii, vkk] * B[vjj, vkk]
```



# Meta Schedule: Auto Tensorization in 4 Steps!

- Step 4: Working with other automatic rules + automatic tensorization

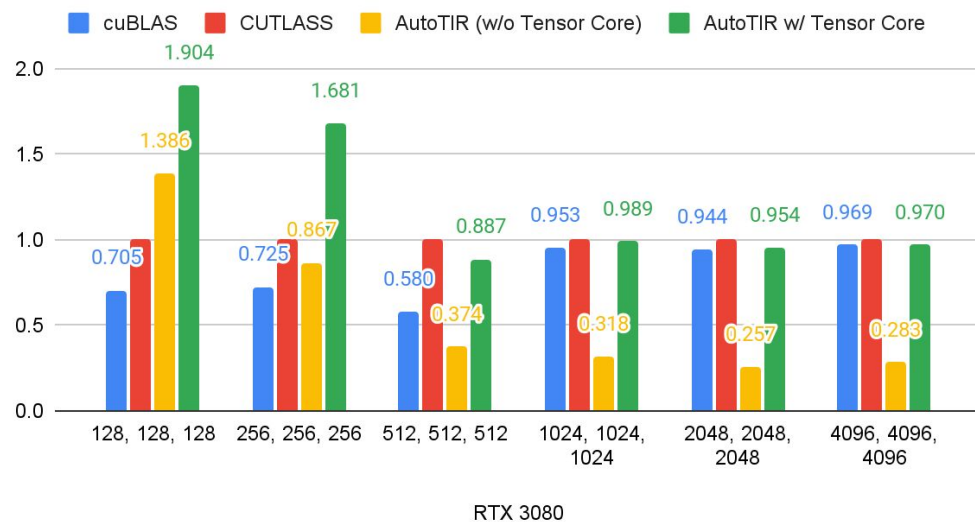
```
@tvm.script.tir
def batch_matmul(a: ty.handle, b: ty.handle, c: ty.handle):
    ...
    for i0, i1_o, i2_o, i3_o in tir.grid(16, 8, 8, 8):
        ...
        for i1_i_init, i2_i_init in tir.grid(16, 16):
            C[vn_init, vi_init, vj_init] = tir.float32(0)
        ...
    tir.evaluate(tir.tvm_mma_sync(
        C.data, tir.floordiv(tir.get_elem_offset(C[vn_1, vi_1, vj_1], dtype="int32"), 256),
        A.data, tir.floordiv(tir.get_elem_offset(A[vn_1, vi_1, vk_1], dtype="int32"), 256),
        B.data, tir.floordiv(tir.get_elem_offset(B[vn_1, vj_1, vk_1], dtype="int32"), 256),
        C.data, tir.floordiv(tir.get_elem_offset(C[vn_1, vi_1, vj_1], dtype="int32"), 256),
        dtype="handle",
    ))
```

# ....and it looks very promising!

(by Bohan from CMU)

- Surpasses cuBLAS perf
- Matches CUTLASS (tuned) perf
- 1.9x on small shapes
- 95% on large shapes
- Tensorized

### GEMM Packed FP16 @ NVIDIA RTX 3080



# Best of both worlds

We care about performance, coverage, and portability, not code generation ideology.



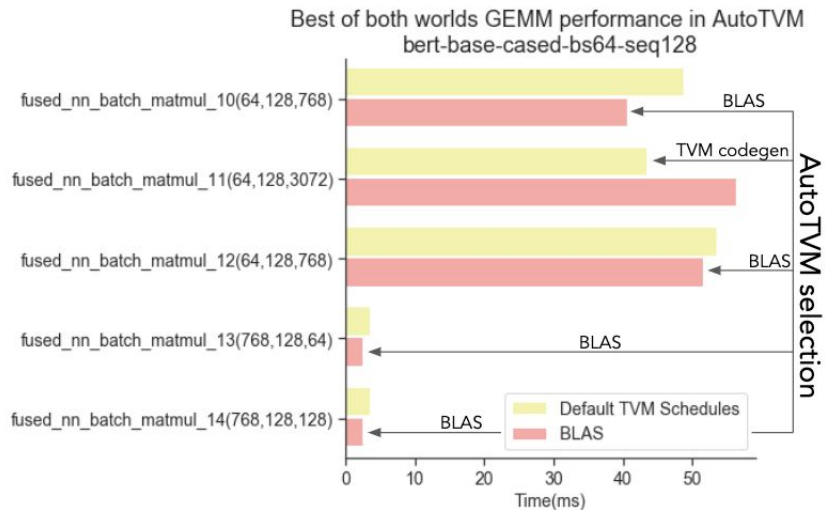
{cuDNN, MKL, ARM  
Compute Lib,  
TensorRT, ...}



# Best of both worlds

For each kernel (or supported subgraph):  
Use `argmax(code generation, kernel library)`

Results: up to 40% gains over TensorRT on Nvidia T4



# Select Performance Results

# Performance at OctoML in 2020/2021

TVM  $\log_2$  fold  
improvement over  
baseline

Over 60 model x hardware benchmarking studies

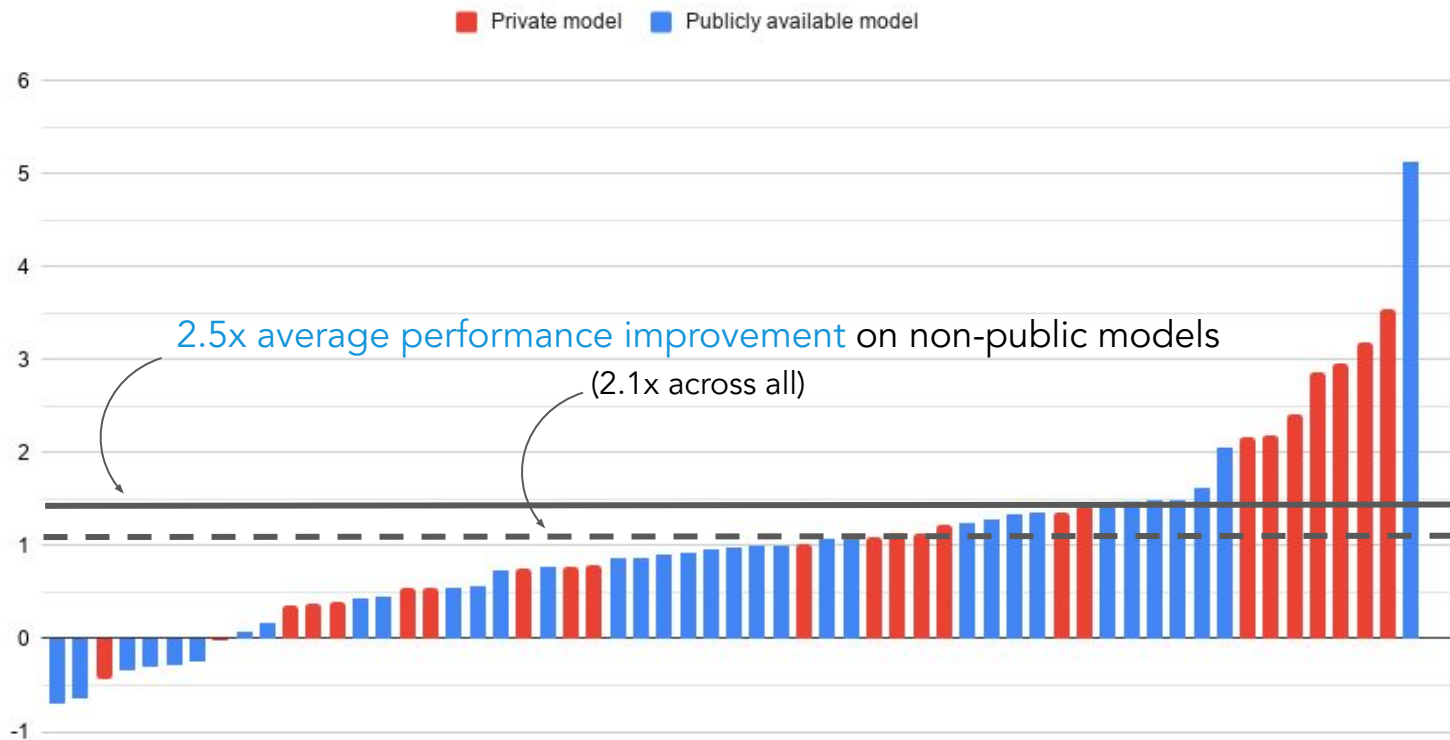
Each study compared TVM against best\* baseline on the target

Sorted by ascending  $\log_2$  gain over baseline

Model x hardware comparison points

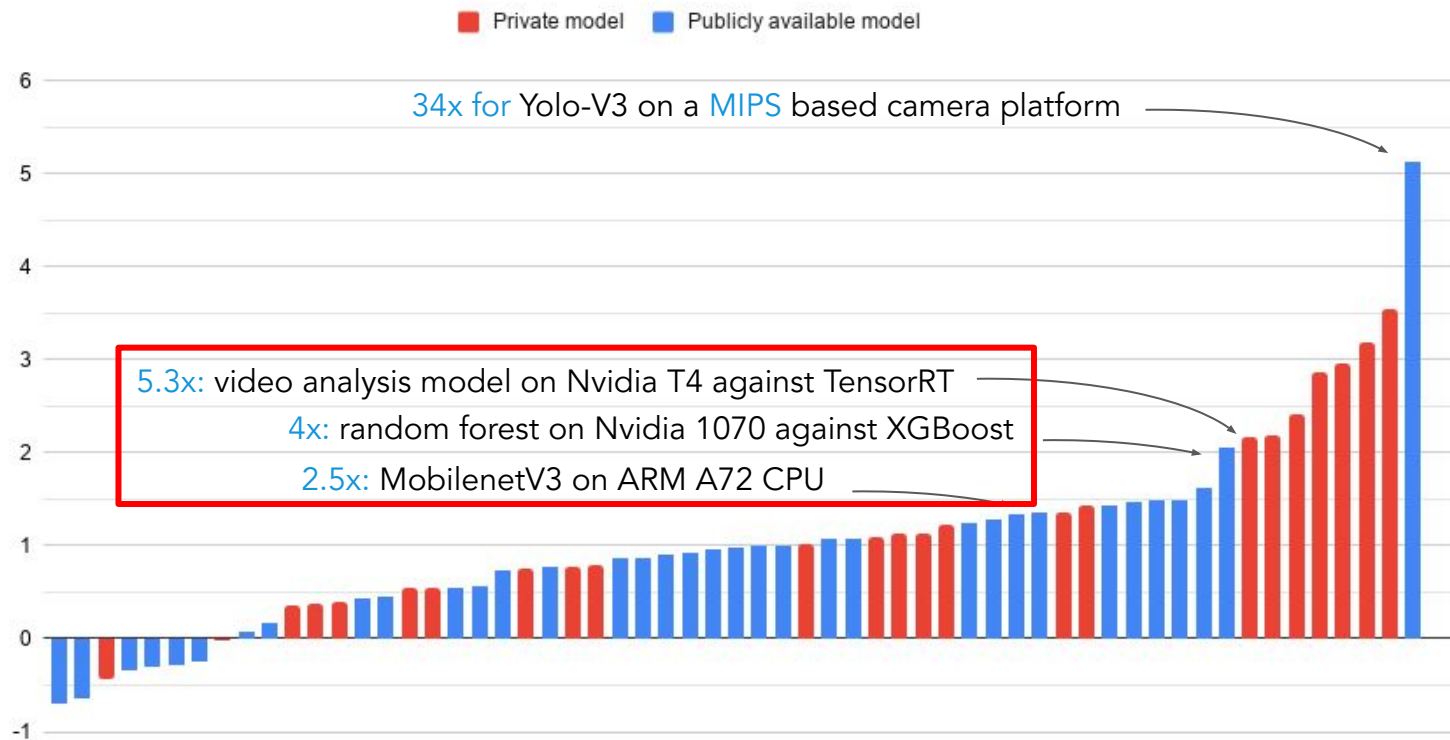
# Higher performance on non-public models

TVM log<sub>2</sub> fold improvement over baseline



Model x hardware comparison points

# Across a broad variety of models and platforms

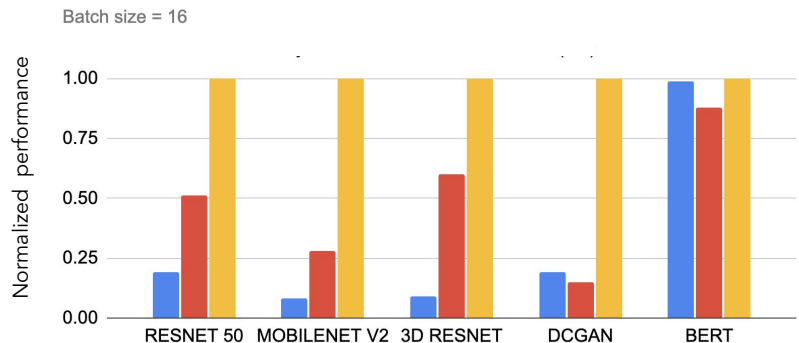
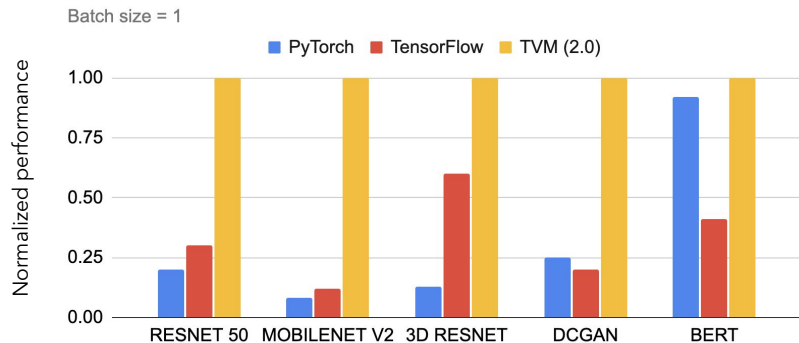


TVM log<sub>2</sub> fold improvement over baseline

Model x hardware comparison points

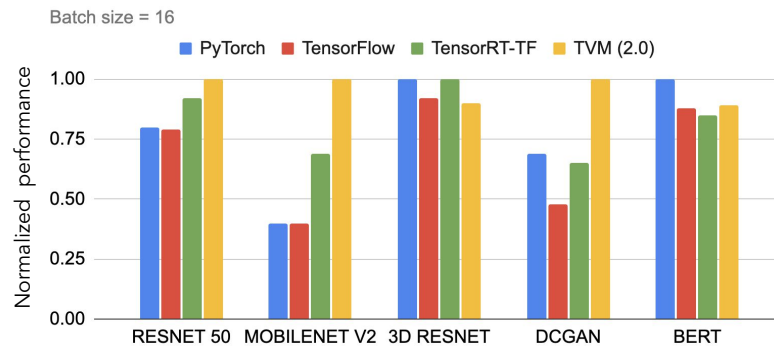
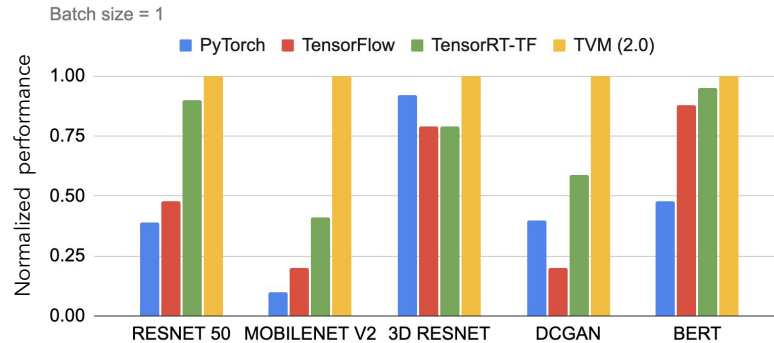
# Results: TVM on popular CPUs and GPUs

## Intel X86 - 2-5X Performance



20 core Intel-Platinum-8269CY fp32 performance data

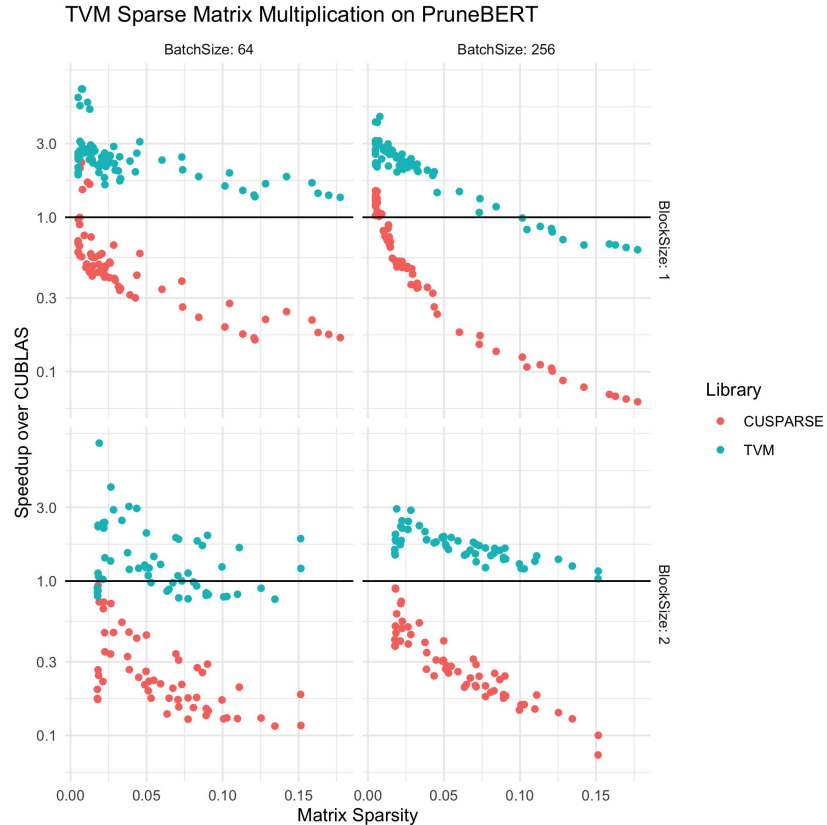
## NVIDIA GPU - 20-50% versus TensorRT



V100 fp32 performance data

# Faster Kernels for Dense-Sparse Multiplication

- Performance comparison on PruneBERT
- 3-10x faster than cuBLAS and cuSPARSE.
- 1 engineer writing TensorIR kernels



# Classical ML (Hummingbird)



<https://github.com/microsoft/hummingbird>

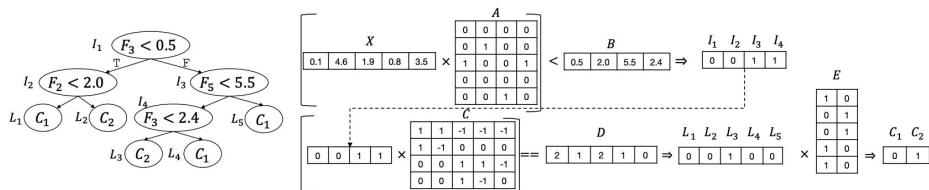
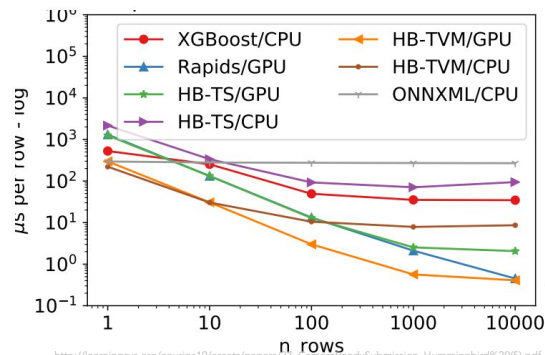


Figure 3: Compiling an example decision tree using the GEMM strategy (algorithm 1).



[http://learningmjs.org/neurips19/assets/papers/27\\_CameraReadySubmission\\_Hummingbird%20\(S\).pdf](http://learningmjs.org/neurips19/assets/papers/27_CameraReadySubmission_Hummingbird%20(S).pdf)

## Tree-models Microbenchmark: Batch w/ GPU

	rf	onnx-ml	hb-pt	hb-ts	hb-onnx	hb-tvm
fraud	2.52s	8.1s	0.15s	0.11s	17.55s	<b>0.02s</b>
year	2.33s	17.23s	0.15s	0.10s	45.12s	<b>0.03s</b>
covtype	47.64s	24.77s	0.32s	0.26s	62.81s	<b>0.06s</b>
epsilon	11.22s	26.03s	0.36s	0.36s	OOM	<b>0.14s</b>
						<b>60x</b>
	xgb	onnx-ml	hb-pt	hb-ts	hb-onnx	hb-tvm
fraud	2.01s	6.4s	0.16s	0.11s	16.89s	<b>0.02s</b>
year	5.77s	15.75s	0.14s	0.1s	44.82s	<b>0.03s</b>
covtype	63.45s	173.92s	1.47s	1.29s	445.89s	<b>0.25s</b>
epsilon	14.84s	29s	0.37s	0.28s	OOM	<b>0.14s</b>
	lgbm	onnx-ml	hb-pt	hb-ts	hb-onnx	hb-tvm
fraud	3.76s	6.59s	0.16s	0.11s	17.70s	<b>0.02s</b>
year	6.18s	10.14s	0.14s	0.10s	OOM	<b>0.03s</b>
covtype	67.12s	158.3s	1.47s	1.29s	446s	<b>0.25s</b>
epsilon	14.13s	26.03s	0.36s	0.28s	OOM	<b>0.14s</b>
						<b>200x</b>

<https://sampl.cs.washington.edu/tvmconf/slides/2019/E13-Matteo-Interlandi.pdf>

## Tree-models Microbenchmark: Batch Inference on CPU

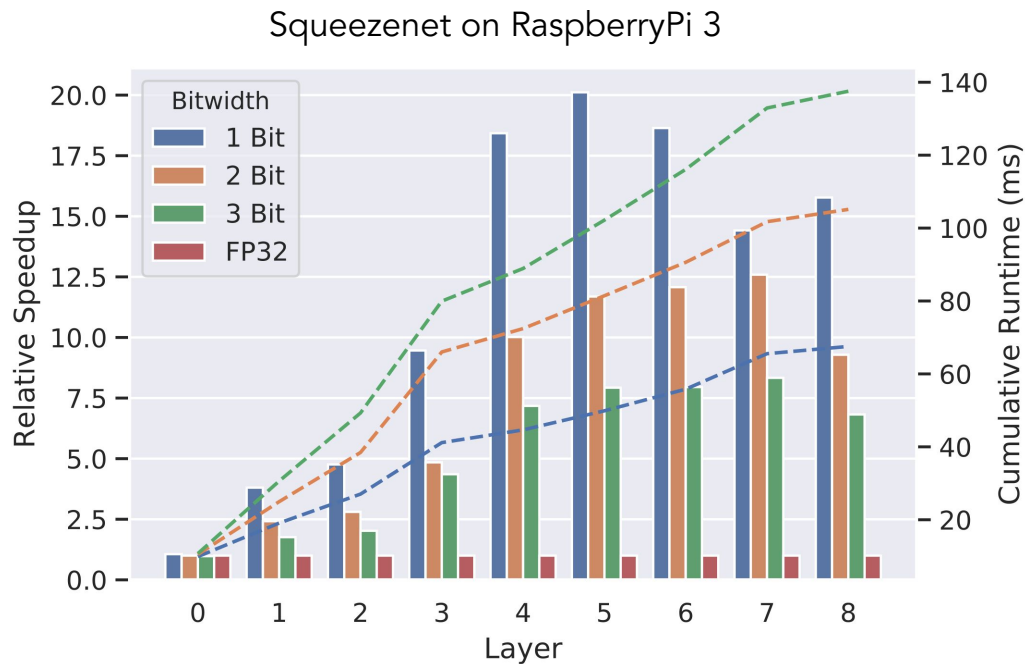
	rf	onnx-ml	hb-pt	hb-ts	hb-onnx	hb-tvm
fraud	2.52s	8.1s	17.18s	17.28	92.58s	<b>3.84s</b>
year	2.33s	17.23s	17.95s	17.23s	154.71s	<b>1.43s</b>
covtype	47.64s	24.77s	38.27s	38.02s	260.55s	<b>20.18s</b>
epsilon	11.22s	26.03s	48.52s	48.87	1255s	<b>8.17s</b>
						<b>50%</b>
	xgb	onnx-ml	hb-pt	hb-ts	hb-onnx	hb-tvm
fraud	2.01s	6.4s	17.23s	16.38s	89.71s	<b>1.93s</b>
year	5.77s	15.75s	17.26s	15.74s	153.96s	<b>1.77s</b>
covtype	63.45s	173.92s	295.6s	295.3s	1255s	<b>28.53s</b>
epsilon	14.84s	29s	47.38s	48.78s	SEGFAULT	<b>4.43s</b>
	lgbm	onnx-ml	hb-pt	hb-ts	hb-onnx	hb-tvm
fraud	3.76s	6.59s	17.41s	16.43s	89.90s	<b>1.97s</b>
year	6.18s	10.14s	18.3s	18.01s	153.67s	<b>1.78s</b>
covtype	67.12s	158.3s	296s	294s	1256s	<b>29.19s</b>
epsilon	14.13s	26.03s	47.21s	47.89s	SEGFAULT	<b>4.41s</b>
						<b>3x</b>

<https://sampl.cs.washington.edu/tvmconf/slides/2019/E13-Matteo-Interlandi.pdf>



# Ultra low bit-width quantization

- In addition to fp32, fp16, int8
- TVM supports bitserial ultra low bit code generation
  - Int{4,3,2,1}
- See Josh Fromm's MLSys 2020 talk and [paper](#).



# Case Study: 90% cloud inference cost reduction

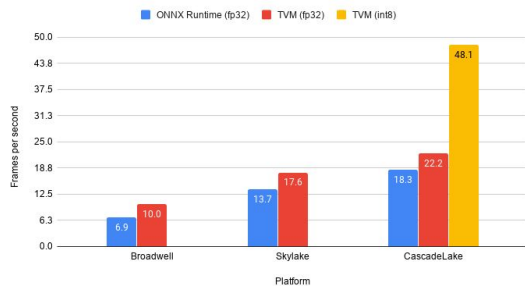
## Background

- Top 10 Tech Company running multiple variations of customized CV models
- Model in batch processing /offline mode using standard HW targets of a major public cloud.
- Billions of inferences per month
- Benchmarking on CPU and GPU

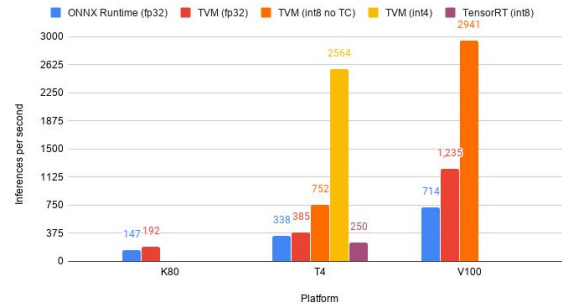
## Results

- 3.8x - TensorRT 8bit to TVM 8bit
- 10x - TensorRT 8bit to TVM 4bit
- **Potential to reduce hourly costs by 90%**

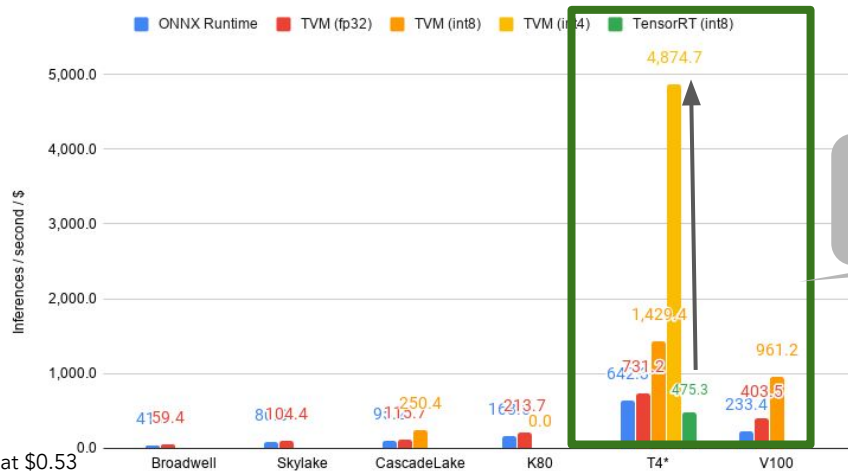
Max throughput - CPU only



Max throughput - GPU only

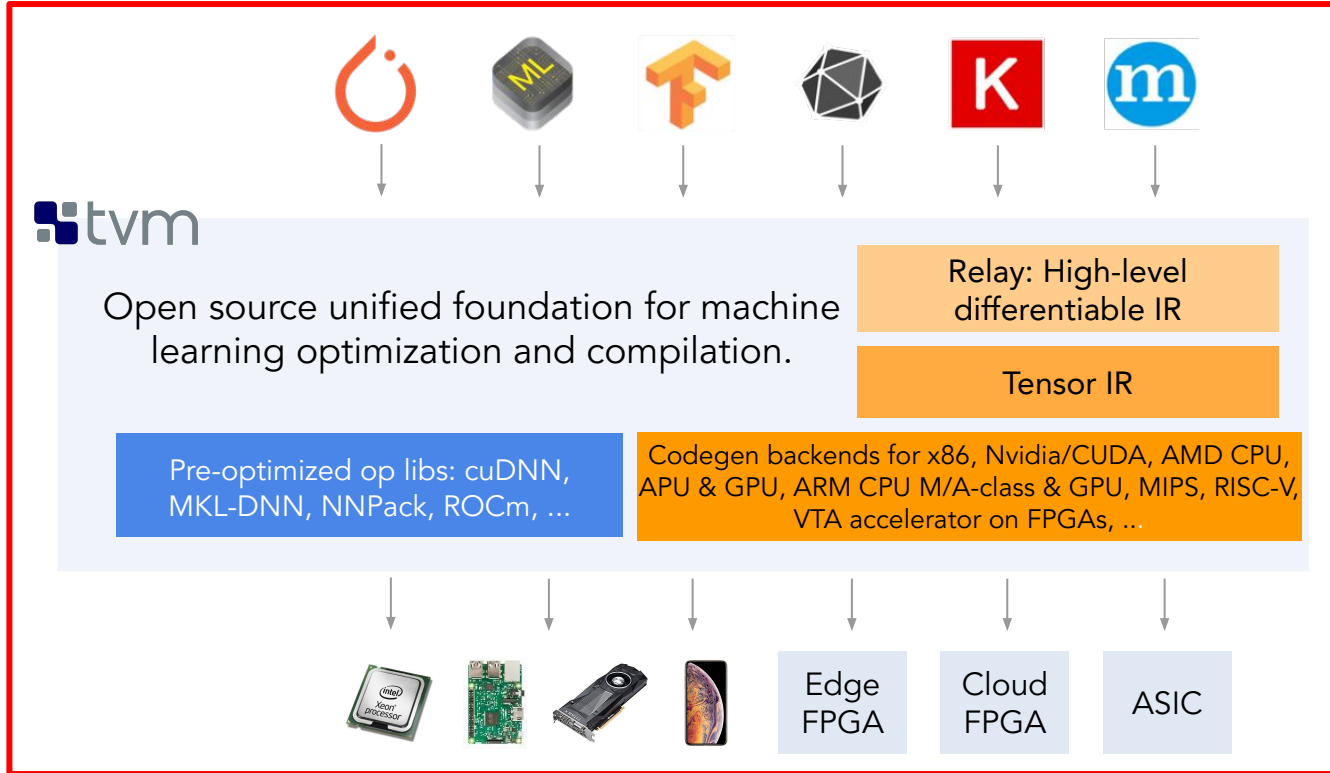


Throughput/\$



Up to 10X  
inferences/dollar  
increase

# Making the most out of TVM & other stacks...



How do end users choose between all the possible configurations of:

- Model definitions
- Optimizations
- Kernel Combinations
- Target Devices
- And more

# Automated platform for ML acceleration & deployment

Models: TensorFlow, PyTorch, ONNX, ...



Hosted service: **No user infrastructure needed.**

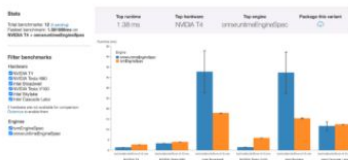
Tuned optimization: **Leveraging OctoML's ML tuning data.**

Automated/flexible packaging: **Easy integration.**

## Optimized Model

- Device-specific tuning for each HW platform
- Leverages tuning data from similar models, HW

## Benchmarking Data



## Packaging

- TVM C Runtime and C API
- Python API
- gRPC
- Docker

# Octomizer Demo

## API access

```
# Specify model file and input layer parameters.
model_file = "mnist.onnx"
input_shapes = {"Input3": (1, 3, 28, 28)}
input_dtypes = {"Input3": "float32"}
model = onnx_model.ONNXModel(client, input_shapes, input_dtypes)

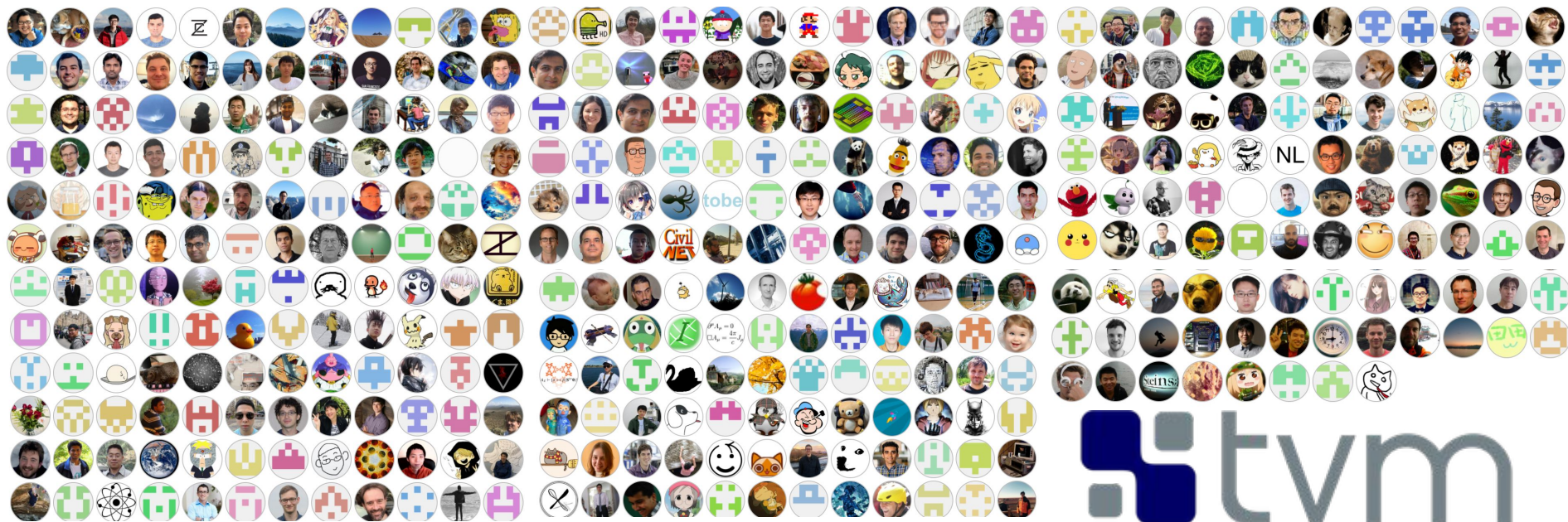
# Upload the file to Octomizer.
modelvar = model.upload("mnist.onnx")

# Optimize it.
wrkflow = modelvar.octomize(platform="broadwell")
wrkflow.wait()
wrkflow.save_package("mnist-octomized.whl")
```

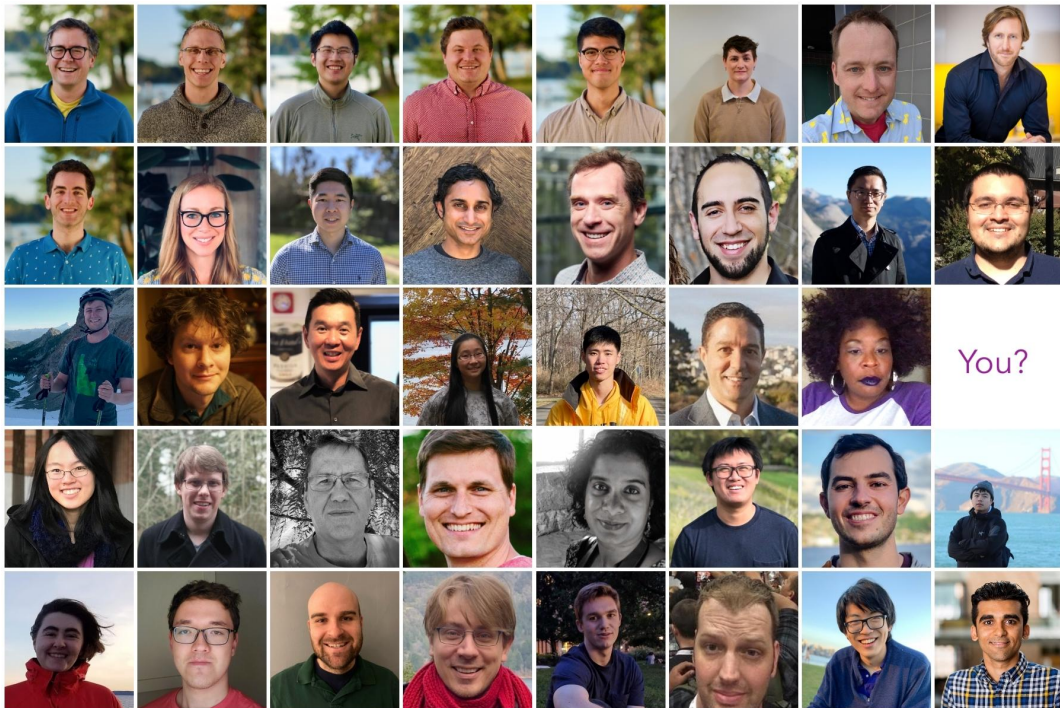
Waitlist! <https://octoml.ai>



# Thank you Apache TVM community! 600+!



# The Octonauts!



You?

<https://octoml.ai/careers>

*Thank you!*

Improving Model Performance, Portability and Productivity  
with Apache TVM and the OctoML platform

*ModSim'21*

*Luis Ceze*

Co-founder & CEO, OctoML.

Professor, University of Washington.



# TVM, ONNX-Runtime, and MLIR

